The uniform general signed rank test and its design sensitivity

Steven R. Howard

Samuel D. Pimentel

April 19, 2019

Abstract

A sensitivity analysis in an observational study tests whether the qualitative conclusions of an analysis would change if we were to allow for the possibility of limited bias due to confounding. The design sensitivity of a hypothesis test quantifies the asymptotic performance of the test in a sensitivity analysis against a particular alternative. We propose a new, non-asymptotic, distribution-free test, the uniform general signed rank test, for observational studies with paired data, and examine its performance under Rosenbaum's sensitivity analysis model. Our test can be viewed as adaptively choosing from among a large underlying family of signed rank tests, and we show that the uniform test achieves design sensitivity equal to the maximum design sensitivity over the underlying family of signed rank tests. Our test thus achieves superior, and sometimes infinite, design sensitivity, indicating it will perform well in sensitivity analyses on large samples. We support this conclusion with simulations and a data example, showing that the advantages of our test extend to moderate sample sizes as well.

1 Introduction

In the empirical study of causal effects, the use of standard statistical hypothesis tests, along with their concomitant *p*-values and confidence intervals, accounts only for the uncertainty introduced by sampling variability. However, in an observational study where treatment assignment has not been randomized, hidden biases due to unobserved confounding can be much larger than sampling uncertainty. As such, standard hypothesis tests may fail to be convincing if they assume the study is free of hidden bias, as a randomized experiment would be. A sensitivity analysis addresses this problem by formally testing whether the qualitative conclusions of a standard procedure would change if hidden bias of a certain magnitude were present (Rosenbaum, 2002).

When an investigator plans to run a sensitivity analysis, the choice of test statistic may no longer hinge solely on traditional measures such as Pitman efficiency. In particular, an investigator may seek a test statistic which is least sensitive to hidden bias, and thereby most likely to successfully distinguish treatment effects from bias, rather than one which is most likely to detect treatment effects in the absence of hidden bias. Design sensitivity is one way to quantify this idea for a particular test statistic (Rosenbaum, 2004, 2010*a*). Design sensitivity complements Pitman efficiency and other conventional means of comparing tests.

Rosenbaum (2010b) shows that a test statistic which focuses on a strongly-affected subgroup may achieve superior design sensitivity, as compared to a statistic which uses all observations. Rosenbaum (2012) shows that a particular test, Noether's test, has excellent design sensitivity but poor power against small effects. Rosenbaum then proposes an adaptive test in which the *p*-value is given by the minimum of two *p*-values from two competing test statistics, correcting for multiple testing by analyzing the joint distribution of these two test statistics. This adaptive test is shown to get some of the best of both worlds, in terms of good power in small samples as well as high design sensitivity. In fact, the adaptive test attains the maximum design sensitivity of its two component tests. Rosenbaum and Small (2017) similarly propose an adaptive test which chooses from the better of two test statistics, one focused on a subgroup and one examining the entire population, with correction for multiple testing.

In this paper we examine a different adaptive test for paired data, in which we may adaptively choose from a large, highly dependent family of test statistics. We control for multiple testing using a uniform concentration bound for the stochastic process formed by this family of test statistics. This permits adaptively choosing among as many test statistics as we have observations, while achieving non-asymptotic, distribution-free

error control. Our theoretical results characterize how this test achieves excellent design sensitivity, which can be infinite against light-tailed alternatives—that is, no matter the strength of confounding, the test will reject with probability approaching one asymptotically. We are not aware of previous discussion of such behavior.

The structure of this paper is as follows. After summarizing Rosenbaum's sensitivity analysis model in Section 2, we describe our test in Section 3, proving that it achieves the promised Type I error control in Theorem 1. We then characterize its design sensitivity with Theorems 2 and 3 of Section 4. Section 5 gives simulation results for a variety of fixed-sample and uniform tests under several light- and heavy-tailed alternatives. We outline the handling of tied data in Section 6, while in Section 7 we illustrate the performance of our tests on an observational dataset examining the link between fish consumption and mercury concentration in the blood. Section 8 concludes and offers some promising avenues for future work.

2 Background and notation

2.1 Rosenbaum's sensitivity analysis model for paired data

We begin with a review of Rosenbaum's sensitivity analysis model for paired data (Rosenbaum, 2002). We have *n* pairs of subjects. The subjects in the *i*th pair have control potential outcomes R_{Cij} , treatment potential outcomes R_{Tij} , and treatment indicators Z_{ij} for j = 1, 2 and $i \in [n]$. Let \mathcal{F} be the σ -field generated by all the potential outcomes $(R_{Cij}, R_{Tij})_{i \in [n], j \in [2]}$.

A sensitivity analysis allows us to test whether a positive conclusion of our study—that is, a rejection of the null—holds up under the possibility of limited confounding. To operationalize this notion, for each $\Gamma \geq 1$ we define the sensitivity analysis null hypothesis $H_0(\Gamma)$, which asserts that

- $R_{Ti1} = R_{Ci1}$ and $R_{Ti2} = R_{Ci2}$ for all $i \in [n]$, i.e., Fisher's sharp null, and
- conditional on \mathcal{F} , treatment assignments are independent between pairs, and the treatment probabilities within each pair are related by the following odds ratio bounds:

$$\frac{1}{\Gamma} \leq \frac{\mathbb{P}\left(Z_{i1}=1 \mid \mathcal{F}\right) / \mathbb{P}\left(Z_{i1}=0 \mid \mathcal{F}\right)}{\mathbb{P}\left(Z_{i2}=1 \mid \mathcal{F}\right) / \mathbb{P}\left(Z_{i2}=0 \mid \mathcal{F}\right)} \leq \Gamma, \quad \text{for all } i \in [n].$$

$$\tag{1}$$

At $\Gamma = 1$, this specifies that, within each pair, both units have the same (conditional) probability of treatment. This is the standard assumption which leads to valid randomization inference in the absence of hidden bias (Rosenbaum, 2002, §3.2).

Write $R_{ij}^{\text{obs}} \coloneqq Z_{ij}R_{Tij} + (1 - Z_{ij})R_{Cij}$ for the observed outcomes and $Y_i = (Z_{i1} - Z_{i2})(R_{i1}^{\text{obs}} - R_{i2}^{\text{obs}})$ for the observed treated-minus-control difference in the *i*th pair. Under $H_0(\Gamma)$ we know that $Y_i = \pm |R_{Ci1} - R_{Ci2}|$ and

$$\frac{1}{1+\Gamma} \le \mathbb{P}\left(Y_i > 0 \mid \mathcal{F}, Z_{i1} + Z_{i2} = 1\right) \le \frac{\Gamma}{1+\Gamma},\tag{2}$$

where for simplicity we assume $\mathbb{P}(Y_i = 0) = 0$ for all *i* throughout this paper. In words, $H_0(\Gamma)$ asserts that there is no effect of treatment for any individual, but the treatment probabilities may differ within a pair in ways we cannot observe. This difference in treatment probabilities could introduce hidden bias into our estimates of the effect of treatment, but the magnitude of such bias is limited by the sensitivity parameter Γ . Again, $\Gamma = 1$ recovers the standard null hypothesis which assumes no hidden bias is present, in which case $\mathbb{P}(Y_i > 0 | \mathcal{F}, Z_{i1} + Z_{i2} = 1) = 1/2$. Throughout the rest of this paper, we implicitly condition on the event $\{Z_{i1} + Z_{i2} = 1, \forall i \in [n]\}$, and omit it from the notation.

This sensitivity analysis model provides a method to conduct valid hypothesis tests under limited confounding, but leaves open the choice of test statistic. In order to judge the relative benefits of different test statistics, we perform a power calculation, comparing the power of various test statistics in a test of the sensitivity analysis null $H_0(\Gamma)$. As with all power calculations, we must choose a particular alternative hypothesis under which to compute power. We define a "favorable" alternative hypothesis $H_1(G)$ for a distribution G over \mathbb{R} , motivated by the following scenario:



Figure 1: The four score functions $\varphi(q)$ used in this paper.

- R_{Cij} is an independent draw from some distribution F, for each $i \in [n], j = 1, 2$,
- $R_{Tij} = R_{Cij} + \tau_i$ for all i, j, where $\tau_i \in \mathbb{R}$ is drawn from some fixed distribution for each $i \in [n]$, and is constant within each pair; and
- $\mathbb{P}(Z_{i1} = 1, Z_{i2} = 0 | \mathcal{F}) = \mathbb{P}(Z_{i1} = 0, Z_{i2} = 1 | \mathcal{F}) = 1/2$, with treatment (conditionally) independent between pairs.

In words, there is a constant treatment effect within pairs and no hidden bias due to unequal treatment probabilities. The alternative hypothesis $H_1(G)$ is then characterized by the induced distribution G of the i.i.d. pair differences $Y_i = (Z_{i1} - Z_{i2})(R_{Ci1} - R_{Ci2}) + \tau_i$; because there is no hidden bias, the mean of this distribution (when the mean exists) is the average treatment effect $\mathbb{E}\tau_i$. In most cases, we consider $\tau_i \equiv \tau$ constant across pairs, so that G is the distribution of $R - R' + \tau$, where R and R' are independent draws from F; this distribution is symmetric about τ . We also consider a "rare effects" model in which τ_i is zero for most pairs and equal to some large value for a small proportion of pairs. In this case, G is a mixture with most mass placed on some distribution symmetric about zero, and the remaining mass on a copy of the distribution shifted to the right.

Rosenbaum's sensitivity analysis model is only one of many possible approaches. For some others, refer to Cornfield et al. (1923/2009); Gilbert et al. (2003); Robins et al. (2000); Yu and Gastwirth (2005). See also Fogarty and Small (2016) for the related problem of sensitivity analysis for multiple outcomes within Rosenbaum's model.

2.2 Sensitivity analysis with general signed rank statistics

Let $(Y_{(i)})$ denote the pair differences (Y_i) ordered by absolute value, so that $|Y_{(1)}| \leq |Y_{(2)}| \leq \cdots \leq |Y_{(n)}|$. A general signed rank statistic has the form

$$T_n = \sum_{i=1}^n \varphi\left(\frac{i}{n+1}\right) \mathbf{1}_{Y_{(i)}>0} \tag{3}$$

for some score function $\varphi : (0,1) \to [0,\infty)$ (Lehmann and Romano, 2005; Rosenbaum, 2010b). The score function allows us to place more or less weight on pairs with larger or smaller observed absolute differences. We will consider four score functions in this paper, all illustrated in Figure 1:

- The sign test uses $\varphi(q) \equiv 1$, so that all pairs contribute equally, regardless of rank. In this case T_n simply counts the number of pairs in which the treated unit had a higher outcome.
- The Wilcoxon signed rank test (WSRT) is equivalent to $\varphi(q) = q$ (Rosenbaum, 2010*a*), so that pairs with larger effects contribute more to the test statistic.
- The normal scores test uses $\varphi(q) = \Phi^{-1}((1+q)/2)$, where Φ^{-1} is the standard normal quantile function, $\mathbb{P}(Z \leq \Phi^{-1}(q)) = q$ when $Z \sim \mathcal{N}(0, 1)$. This score function is the quantile function of the absolute value of a standard normal random variable, and this general signed rank test has high power when outcomes are drawn from a normal distribution (Lehmann and Romano, 2005, §6.9-6.10).

• Finally, we include a "**redescending**" score function, $\varphi(q) = \sum_{l=\underline{m}}^{\overline{m}} \frac{l}{m} {m \choose l} q^{l-1} (1-q)^{m-l}$, so-called because this function rises as q increases from zero, like the WSRT and normal scores functions do, but falls back to zero as q approaches one, unlike the other three score functions. The resulting statistic puts more weight on pairs with larger absolute differences, but excludes the most extreme observations, which may be outliers. We set $(m, \underline{m}, \overline{m}) = (20, 12, 19)$. This score function approximates the *U*-statistic described in Rosenbaum (2011, Lemma 1), and the given values of $(m, \underline{m}, \overline{m})$ were found to perform well in Rosenbaum's study.

The sensitivity analysis null hypothesis $H_0(\Gamma)$ does not specify a single distribution for the observables (Y_i) , but it does imply a single worst-case distribution for the test statistic T_n in a one-sided test which rejects for T_n sufficiently large—that is, a distribution which maximizes $\mathbb{P}(T_n \ge a \mid \mathcal{F})$ for any threshold a, among all distributions in $H_0(\Gamma)$. This worst-case distribution has the n signs $(1_{Y_i>0})$ independent with $\mathbb{P}(Y_i > 0 \mid \mathcal{F}) = \Gamma/(1 + \Gamma)$ for all $i \in [n]$ (Rosenbaum, 2002, §4.3). Write $c_{\alpha,n}(\Gamma)$ for the $1 - \alpha$ quantile of T_n under this worst-case distribution, so that $c_{\alpha,n}(\Gamma)$ is the critical value of a one-sided, level- α sensitivity analysis testing $H_0(\Gamma)$ with test statistic T_n ; the critical value may depend on \mathcal{F} , in the case of ties. This critical value yields a valid (conditional) test of the sensitivity analysis null hypothesis, and is not hard to approximate numerically or via the normal distribution. In Theorem 1 below, we build upon these ideas to define a uniform general signed rank test, deriving closed-form critical values which guarantee non-asymptotic Type I error control under the sensitivity null $H_0(\Gamma)$.

2.3 Power of a sensitivity analysis and design sensitivity

Under $H_1(G)$, the power of a one-sided, level- α sensitivity analysis for a general signed rank test with statistic T_n is $\mathbb{P}_1(T_n \geq c_{\alpha,n}(\Gamma))$, which is well-defined since $H_1(G)$ specifies the distribution of T_n completely. This power depends on the level α , the sample size n, the sensitivity parameter Γ , the alternative distribution G, and the score function φ . The design sensitivity (Rosenbaum, 2004, 2010a) of the test statistic T_n is the value $\widetilde{\Gamma}$ such that, as the sample size grows without bound, the power of a sensitivity analysis with parameter Γ approaches one whenever $\Gamma < \widetilde{\Gamma}$ and approaches zero whenever $\Gamma > \widetilde{\Gamma}$:

$$\lim_{n \to \infty} \mathbb{P}_1(T_n \ge c_{\alpha,n}(\Gamma)) = 1, \quad \text{for} \quad 1 \le \Gamma < \widetilde{\Gamma}, \quad \text{and}$$
(4)

$$\lim_{n \to \infty} \mathbb{P}_1(T_n \ge c_{\alpha,n}(\Gamma)) = 0, \quad \text{for} \quad \widetilde{\Gamma} < \Gamma < \infty.$$
(5)

Formally, the design sensitivity depends on the level α , the alternative distribution G and the score function φ . In typical examples, including those considered below, the dependence on α vanishes. It is clear from the definition that such a value is unique, if it exists, but existence must be proved as part of the derivation of design sensitivity, as in our Theorem 2. Note also that we may have $\tilde{\Gamma} = \infty$, which means that $\lim_{n\to\infty} \mathbb{P}_1(T_n \ge c_{\alpha,n}(\Gamma)) = 1$ for all $\Gamma \ge 1$; in words, the test has power approaching one against the given alternative regardless of how large a sensitivity parameter Γ is chosen.

Proposition 2 of Rosenbaum (2010b) gives a formula for the design sensitivity of a general signed rank test whenever the score function φ is piecewise continuous, nondecreasing and not identically zero:

$$\widetilde{\Gamma} = \frac{\pi}{1 - \pi}, \quad \text{where} \quad \pi \coloneqq \frac{\int_0^\infty \varphi(G(y) - G(-y)) \, \mathrm{d}G(y)}{\int_0^1 \varphi(y) \, \mathrm{d}y}.$$
(6)

Note that G(y) - G(-y) is the CDF of |Y| under $H_1(G)$. We see that the design sensitivity of a general signed rank test is determined precisely by the aspects of φ and G captured in the quantity π . In Theorems 2 and 3, we extend this result to characterize the design sensitivity of our uniform general signed rank test. Our conditions on φ , while not strictly more general, do allow for the normal scores and redescending score functions, in contrast to Rosenbaum's conditions.

For the sign test, $\varphi(q) \equiv 1$, we have $\int_0^1 \varphi(y) \, dy = 1$ and $\int_0^\infty \varphi(G(y) - G(-y)) \, dG(y)$ is exactly $\mathbb{P}(Y > 0)$ when $Y \sim G$. Hence $\pi = \mathbb{P}_1(Y > 0)$ (cf. Rosenbaum, 2012, Proposition 1). In words, this π is simply the probability that a pair difference Y gives evidence in favor of a positive treatment effect, under the favorable alternative with no hidden bias.

3 A uniform general signed rank test

We now define a general class of uniform signed rank tests which operate on a family of related test statistics $(T_n(x))_{x \in (0,1)}$. Informally, our test rejects when *any* test statistic in the family lies above a corresponding modified critical value. These critical values are carefully chosen to correct for multiplicity by taking advantage of the structure of the family of test statistics. The uniform nature of our test yields advantages in terms of design sensitivity, which we describe in Section 4.

For any $\varphi : (0,1) \to [0,\infty)$, define the family of test statistics $(T_n(x))_{x \in (0,1)}$ by $T_n(x) = 0$ for x < 1/(n+1), and for $x \ge 1/(n+1)$,

$$T_n(x) \coloneqq \sum_{i=\lceil (1-x)(n+1)\rceil}^n \varphi\left(\frac{i}{n+1}\right) \mathbf{1}_{Y_{(i)}>0} = \sum_{i=\lceil (1-x)(n+1)\rceil}^n c_i \mathbf{1}_{Y_{(i)}>0},\tag{7}$$

where we have defined $c_i \coloneqq \varphi\left(\frac{i}{n+1}\right)$ for convenience. For each x, $T_n(x)$ is a general signed rank statistic using the "truncated" score function $\varphi_x(q) = \varphi(q) \mathbf{1}_{q \ge 1-x}$. There are n distinct nontrivial test statistics in this family, $T_n(k/(n+1))$ for $k = 1, \ldots, n$, corresponding to the partial sums $\sum_{i=k}^n c_i \mathbf{1}_{Y_{(i)}>0}$ for $k = n, n-1, \ldots, 1$. Hence the family corresponds to a random walk with n steps and step sizes determined by the function $\varphi(\cdot)$.

Note that, despite the generality of our construction in terms of the score function φ , our family always consists of truncated versions of the full test statistic. Such truncated statistics focus on subsets of the experimental sample with large observed effects $|Y_i|$. As such, our test will tend to perform especially well against alternatives with large, rare effects.

Our uniform test will be characterized by a threshold function $f_{\alpha,n}(x)$, the uniform analogue of a critical value. Our test rejects whenever $T_n(x) \ge f_{\alpha,n}(x)$ for any $x \in (0,1)$. As in the fixed-sample case, there is a single worst-case distribution under $H_0(\Gamma)$ which maximizes the probability of rejection; we prove the following in Appendix A.1.

Proposition 1. Fix any threshold function $f_{\alpha,n} : (0,1) \to \mathbb{R}_{>0}$. Among all distributions in $H_0(\Gamma)$, the rejection probability $\mathbb{P}(\exists x \in (0,1) : T_n(x) \ge f_{\alpha,n}(x) \mid \mathcal{F})$ is maximized when $\mathbb{P}(Y_i > 0 \mid \mathcal{F}) = \Gamma/(1+\Gamma)$ for all $i \in [n]$.

Under this worst-case distribution in $H_0(\Gamma)$, each step of the random walk equals c_i with probability $\rho_{\Gamma} := \Gamma/(1+\Gamma)$ and zero otherwise; these steps are independent. The resulting mean and variance of $T_n(x)$ are

$$\mu_n(x) \coloneqq \mathbb{E}T_n(x) = \rho_\Gamma \sum_{i=\lceil (1-x)(n+1)\rceil}^n c_i \tag{8}$$

$$\sigma_n^2(x) \coloneqq \operatorname{Var} T_n(x) = \rho_{\Gamma}(1 - \rho_{\Gamma}) \sum_{i = \lceil (1 - x)(n + 1) \rceil}^n c_i^2.$$
(9)

Our threshold function requires a tuning parameter $x_0 > 0$ to be fixed in advance, such that $\sigma_n^2(x_0) > 0$. If $\sigma_n^2(x) = 0$ for all x, then we cannot choose a valid x_0 , but in this case, $T_n(x) = 0$ a.s. for all x, so we cannot reject for any reasonable bound. We then construct the following high-probability uniform upper boundary on the random walk $T_n(x)$:

$$f_{\alpha,n}(x) \coloneqq \frac{1}{\lambda_n} \left[\log\left(\frac{1}{\alpha}\right) + \sum_{i=\lceil (1-x)(n+1)\rceil}^n \log\left(1 + \rho_{\Gamma}(e^{c_i\lambda_n} - 1)\right) \right], \quad \text{where } \lambda_n \coloneqq \sqrt{\frac{2\log\alpha^{-1}}{\sigma_n^2(x_0)}}. \tag{10}$$

For notational simplicity, we omit the dependence of $f_{\alpha,n}$ on x_0 .

Theorem 1. Under $H_0(\Gamma)$, for any $x_0 > 0$ such that $\sigma_n^2(x_0) > 0$ and any $\alpha \in (0,1)$, we have

$$\mathbb{P}\left(\exists x \in (0,1) : T_n(x) \ge f_{\alpha,n}(x) \mid \mathcal{F}\right) \le \alpha.$$
(11)

Theorem 1 justifies rejecting the sensitivity null $H_0(\Gamma)$ whenever $T_n(x) \ge f_{\alpha,n}(x)$ for some $x \in (0, 1)$, allowing us to adaptively choose a value of x after seeing the data, while retaining Type I error control at level α . We



Figure 2: Illustration of Theorem 1 and the uniform bound (10) for the uniform sign test, $\varphi(q) \equiv 1$. Black line shows one realization of the random walk $T_n(x)$ for $x = 1/(n+1), 2/(n+1), \ldots, n/(n+1)$; here n = 50 and $\Gamma = 2$. Green line shows the uniform upper bound $f_{\alpha,n}(x)$ which is unlikely to ever be crossed by the random walk. We may think of each value $f(1/(n+1)), f(2/(n+1)), \ldots$ as a modified critical value for the corresponding test statistic.

call this test a uniform general signed rank test. The idea is illustrated in Figure 2. Because the probability bound in Theorem 1 holds uniformly over all x, in any given dataset we may choose the value of x which yields the strongest inference. We can think of the resulting test as simultaneously conducting general signed rank tests with truncated score functions $\varphi_x(q) = \varphi(q) \mathbf{1}_{q \ge 1-x}$ for all values $x = 1/(n+1), \dots, n/(n+1)$, but with modified critical values given by $f_{\alpha,n}(x)$. The critical value $f_{\alpha,n}(x)$ is larger than the fixed-sample exact critical value $c_{\alpha,n}(\Gamma)$ from Section 2.2, accounting for the uniformity of our test. Note that, when we use the sign test score function $\varphi(q) = 1$, the resulting truncated score functions φ_x are exactly the score functions used in Noether's test (Noether, 1973; Rosenbaum, 2012).

Before proving Theorem 1 we give some intuition for the bound $f_{\alpha,n}$ based on the following asymptotic approximation, which holds under mild conditions on φ as detailed in Appendix A.3:

$$f_{\alpha,n}(x) = \underbrace{\mu_n(x) + \left(1 + \frac{\sigma_n^2(x)}{\sigma_n^2(x_0)}\right) \sqrt{\frac{\sigma_n^2(x_0) \log \alpha^{-1}}{2}}_{g_{\alpha,n}(x)} + \mathcal{O}(1).$$
(12)

The leading term, $\mu_n(x)$, is $\mathcal{O}(n)$ and accounts for the drift of the random walk. The next term is $\mathcal{O}(\sqrt{n})$ and accounts for the deviations of the random walk about its mean. As discussed in Appendix A.3, the parameter x_0 determines the value of x for which the boundary $g_{\alpha,n}(x)$ is optimized, and this motivates the choice of λ_n in the definition of $f_{\alpha,n}$. Theorem 1 would continue to hold with any choice $\lambda_n > 0$, but our choice yields the interpretable tuning parameter x_0 .

The discussion in Appendix A.3 also shows that the remainder $g_{\alpha,n}(x) - f_{\alpha,n}(x)$ is always negative, so that $g_{\alpha,n}(x)$ yields an alternative threshold function with a simpler analytical form, but the resulting test has slightly less power. In fact, the uniform boundaries $f_{\alpha,n}$ and $g_{\alpha,n}$ are drawn from a broader framework for uniform concentration of random walks described in Howard et al. (2018*a*,*b*). Other boundaries are possible and will yield different performance; further exploration of alternative boundaries is a promising avenue of future work. We give below a short, self-contained proof of Theorem 1 to illustrate the techniques, which are closely related to the classical Cramér-Chernoff method (Cramér, 1938; Chernoff, 1952; Boucheron et al., 2013, section 2.2).

Proof of Theorem 1. Throughout the proof, we condition on \mathcal{F} , dropping it from the notation for simplicity. Let $S_i \coloneqq 1_{Y_{(i)}>0}$ for $i \in [n]$, so that $T_n(k/(n+1)) = \sum_{i=n+1-k}^n c_i S_i$ for each $k \in [n]$. By Proposition 1, under the worst-case distribution in $H_0(\Gamma)$, $(S_i)_{i \in [n]}$ are distributed as n i.i.d. Bernoulli (ρ_{Γ}) random variables. The moment-generating function of the random variable $c_i S_i$ is

$$\mathbb{E}e^{\lambda c_i S_i} = 1 + \rho_{\Gamma}(e^{c_i \lambda} - 1) \quad \forall \lambda \in \mathbb{R}.$$
(13)

Now define $(L_k)_{k=0}^n$ by $L_0 \coloneqq 1$ and, for $k \in [n]$,

$$L_k \coloneqq \exp\left\{\lambda_n T_n\left(\frac{k}{n+1}\right) - \sum_{i=n+1-k}^n \log\left(1 + \rho_{\Gamma}(e^{c_i\lambda_n} - 1)\right)\right\} = \prod_{i=n+1-k}^n \frac{e^{\lambda_n c_i S_i}}{1 + \rho_{\Gamma}(e^{c_i\lambda_n} - 1)}.$$
 (14)

It is easy to see from (13) and (14) that $\mathbb{E}(L_k | S_n, S_{n-1}, \dots, S_{n+2-k}) = L_{k-1}$, so that L_k is a nonnegative martingale with respect to the natural filtration defined by the sequence S_n, S_{n-1}, \dots, S_1 . Then Ville's maximal inequality for nonnegative supermartingales (Ville, 1939; Durrett, 2013, Exercise 5.7.1) implies

$$\alpha \ge \mathbb{P}\left(\exists k \in [n] : L_k \ge \alpha^{-1}\right) \tag{15}$$

$$= \mathbb{P}\left(\exists k \in [n] : T_n\left(\frac{k}{n+1}\right) \ge f_{\alpha,n}\left(\frac{k}{n+1}\right)\right)$$
(16)

$$= \mathbb{P}\left(\exists x \in \left\{\frac{1}{n+1}, \frac{2}{n+1}, \dots, \frac{n}{n+1}\right\} : T_n(x) \ge f_{\alpha,n}(x)\right)$$
(17)

$$= \mathbb{P}\left(\exists x \in (0,1) : T_n(x) \ge f_{\alpha,n}(x)\right).$$
(18)

The final equality follows since the values $x = 1/(n+1), 2/(n+1), \ldots, n/(n+1)$ capture all of the distinct values of both $T_n(x)$ and $f_{\alpha,n}(x)$ for $x \ge 1/(n+1)$, and adding the region 0 < x < 1/(n+1) does not change the overall probability since $T_n(x) = 0$ over this region while $f_{\alpha,n}(x)$ is strictly positive. \Box

4 Design sensitivity of the uniform test

=

We have shown that the uniform test may be thought of as simultaneously conducting general signed rank tests at all values of x with modified critical values $f_{\alpha,n}(x)$. We might equivalently think of this as adjusting the significance level α downwards, and to different values for different x, in computing critical values for a sequence of general signed rank tests. Recalling that the design sensitivity of a general signed rank test (6) does not depend on α , we may wonder if the uniform test has design sensitivity equal to the maximum of the design sensitivities of the component test statistics $T_n(x)$. This conclusion is not quite trivial, since the "adjusted significance levels" in the uniform test vary as n grows. Nonetheless, it turns out to be true. We prove this for score functions $\varphi : (0, 1) \to [0, \infty)$ satisfying the following properties:

- (P1) $\int_0^1 \varphi^2(x) \, \mathrm{d}x < \infty;$
- (P2) φ is discontinuous on a set of Lebesgue measure zero;
- (P3) there exists a constant $a \in [0, 1/2)$ such that φ is nonincreasing on (0, a), nondecreasing on (1 a, 1), and bounded on (a, 1 a); and
- (P4) $\int_{1-x}^{1} \varphi(x) \, \mathrm{d}x > 0$ for all x > 0.

Theorem 2. Suppose φ satisfies conditions (P1-P4) above, and G is continuous. Then the design sensitivity of the corresponding uniform general signed rank test under $H_1(G)$ is

$$\widetilde{\Gamma}_{\varphi,\text{unif}} \coloneqq \sup_{x \in (0,1)} \widetilde{\Gamma}(x) = \sup_{x \in (0,1)} \frac{\pi(x)}{1 - \pi(x)}, \quad where \quad \pi(x) \coloneqq \frac{\int_0^\infty \varphi(G(y) - G(-y)) \mathbf{1}_{G(y) - G(-y) \ge 1 - x} \, \mathrm{d}G(y)}{\int_{1 - x}^1 \varphi(y) \, \mathrm{d}y}.$$
(19)

Most of the work in the proof of Theorem 2 is captured by the following pair of lemmas. The first, proved in Appendix A.3, characterizes the asymptotic behavior of the boundary $f_{\alpha,n}(x)$ as $n \to \infty$.

Lemma 1. If φ satisfies conditions (P1)-(P3) above, then for any $x_0 > 0$ such that $\sigma_n^2(x_0) > 0$, any $\alpha \in (0,1)$, and any $x \in (0,1)$, we have $n^{-1}\mu_n(x) \to \rho_{\Gamma} \int_{1-x}^1 \varphi(y) \, dy$ and $f_{\alpha,n}(x) = \mu_n(x) + \mathcal{O}(\sqrt{n})$ as $n \to \infty$.

The second lemma generalizes a result of Sen (1970); we give the proof in Appendix A.4.

Lemma 2. If φ satisfies conditions (P1-P3) above, and Y_1, Y_2, \ldots are drawn i.i.d. from a continuous distribution G, then

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \varphi\left(\frac{i}{n+1}\right) \mathbf{1}_{Y_{(i)} > 0} = \int_{0}^{\infty} \varphi(G(y) - G(-y)) \,\mathrm{d}G(y) \quad a.s.$$
(20)

Proof of Theorem 2. Let $H(x) \coloneqq G(x) - G(-x)$ denote the distribution of |Y|. Fix any $x \in (0, 1)$. Applying Lemma 2 to the truncated score function $\varphi_x(q) = \varphi(q) \mathbf{1}_{q \ge 1-x}$ yields

$$\lim_{n \to \infty} \frac{T_n(x)}{n} = \int_0^\infty \varphi(H(y)) \mathbf{1}_{H(y) \ge 1-x} \,\mathrm{d}G(y) \quad \text{a.s.}$$
(21)

Meanwhile, Lemma 1 implies that

$$\lim_{n \to \infty} \frac{f_{\alpha,n}(x)}{n} = \rho_{\Gamma} \int_{1-x}^{1} \varphi(y) \,\mathrm{d}y.$$
(22)

Combining (22) with (21), we conclude that

$$\mathbb{P}(T_n(x) \ge f_{\alpha,n}(x)) = \mathbb{P}(n^{-1}T_n(x) \ge n^{-1}f_{\alpha,n}(x)) \to 1$$

if
$$\int_0^\infty \varphi(H(y)) \mathbf{1}_{H(y)\ge 1-x} \,\mathrm{d}G(y) > \rho_\Gamma \int_{1-x}^1 \varphi(y) \,\mathrm{d}y, \quad (23)$$

that is, if $\Gamma < \pi(x)/[1-\pi(x)]$. Since the uniform test rejects whenever $T_n(x) \ge f_{\alpha,n}(x)$ for some x, it will reject with probability approaching one whenever $\Gamma < \pi(x)/[1-\pi(x)]$ for some $x \in (0,1)$. By a similar argument, $\mathbb{P}(T_n(x) \ge f_{\alpha,n}(x)) \to 0$ if $\Gamma > \pi(x)/[1-\pi(x)]$, so the uniform test will reject with probability approaching zero if $\Gamma > \pi(x)/[1-\pi(x)]$ for all $x \in (0,1)$. The conclusion follows.

Compare Theorem 2 to Proposition 1 of Rosenbaum (2012). Rosenbaum constructs an adaptive test choosing between two test statistics and achieving design sensitivity equal to the maximum of the two component tests. Theorem 2 shows that this principle may be extended to an infinite family of tests, in this case because the family possesses a dependence structure that allows us to construct an appropriate uniform bound.

We note that all of the score functions introduced in Section 2.2 satisfy conditions (P1-P4). Most of these are obvious; the only work required is to show that the score function for the normal scores test satisfies property (P1), and we give the short proof in Appendix A.5.

Proposition 2. For the normal scores function, $\varphi(q) = \Phi^{-1}((1+q)/2)$, we have $\int_0^1 \varphi^p(x) dx < \infty$ for all $p \ge 1$.

Figure 3 shows $\pi(x)$ as defined in Theorem 2. Each panel includes three alternative distributions G: normal with unit variance, Laplace (double exponential) with unit scale, and Cauchy with unit scale. In the first two panels, each distribution is centered at $\tau = 1/2$. The bottom panel shows a "rare effects" model in which G is a mixture of two of the given base distributions, one centered at zero receiving 90% of the total mass, and the other centered at $\tau = 5$ receiving 10% of the total mass. This simulates a situation in which 90% of pairs have no treatment effect, while the remaining 10% of pairs have a large constant treatment effect, so that the average treatment effect remains equal to 1/2.

The first two panels of Figure 3 show $\pi(x)$ for the sign and WSRT score functions introduced in Section 2.2; Appendix A.7 includes $\pi(x)$ plots for the normal scores and redescending score functions, which are qualitatively similar to $\pi(x)$ for the WSRT. For the sign test, $\pi(x)$ is maximized at some value x < 1 under all distributions, although the increase is modest for the Laplace and Cauchy alternatives. This illustrates the benefits of truncation with the sign test. With the WSRT, we still see dramatic gains under a normal alternative, and indeed $\pi(x) \uparrow 1$ as $x \downarrow 0$ for all of our score functions under a normal alternative. This indicates we can achieve infinite design sensitivity under normal tails, a fact which we prove in Corollary 1. Under the Laplace or Cauchy alternatives, however, we do not see substantial gains in $\pi(x)$ as x decreases from one for the WSRT; the same holds true for the normal scores and redescending score functions. Under the heavier-tailed Laplace and Cauchy alternatives, it seems, score functions which place more weight on larger outcomes do not benefit from narrowing attention to a subset of pairs with the largest absolute differences. Informally speaking, the higher likelihood of large outliers means less information is present in the tails.

The $\pi(x)$ functions in the bottom panel, computed under a rare effects model, tells a different story. Here, a uniform WSRT benefits from narrowing attention to a subset of pairs with large absolute differences regardless of the alternative distribution, although gains are still more modest for the Cauchy alternative than for the others. This confirms the intuitive fact that, when effects are large and rare, a test which restricts attention accordingly yields lower sensitivity to hidden bias.



Figure 3: $\pi(x)$ from Theorem 2 for sign and WSRT score functions when G is standard normal, Laplace (double exponential) or Cauchy. First two panels show alternative with $\tau = 1/2$. Bottom panel shows rare effects model: 90% of pairs have no treatment effect, $\tau = 0$ while 10% of pairs have a large treatment effect, $\tau = 5$. See Figure 6 for corresponding plots with normal scores and redescending score functions, which have $\pi(x)$ qualitatively similar to that for the WSRT score function.

Figure 3 makes it clear that the best choice of x depends on the alternative distribution G and the score function in a complicated manner. The advantage of our uniform test is that it can adapt to the alternative at hand without prior knowledge, achieving performance equivalent to the oracle choice of x in terms of design sensitivity. It it also notable that all four score functions exhibit identical behavior near x = 0. The following result makes this observation precise whenever G is continuous with infinite support. We show that the limiting behavior of $\pi(x)$ as $x \downarrow 0$ is often determined by the tails of G alone, not by the score function φ , and this may be used to lower bound the design sensitivity over a broad class of score functions.

Theorem 3. Suppose φ satisfies conditions (P1-P4) above, and suppose G has positive density g(x) with respect to Lebesgue measure for all $x \in \mathbb{R}$. Then

$$\widetilde{\Gamma}_{\varphi,\mathrm{unif}} \ge \liminf_{q \uparrow \infty} \frac{g(q)}{g(-q)}.$$
(24)

Proof. Write q_x for the x-quantile of |Y| when $Y \sim G$, so that q_x is defined by the equation $G(q_x) - G(-q_x) = x$. We shall require the derivative of q_x below, which we find by implicit differentiation:

$$\frac{\mathrm{d}q_x}{\mathrm{d}x} = \frac{1}{g(q_x) + g(-q_x)}.\tag{25}$$

Now observe that, using the definition of q_x , we may write $\pi(x)$ from Theorem 2 as

$$\pi(x) = \frac{\int_{q_{1-x}}^{\infty} \varphi(G(y) - G(-y)) \,\mathrm{d}G(y)}{\int_{1-x}^{1} \varphi(y) \,\mathrm{d}y}.$$
(26)

We apply the generalized form of L'Hôpital's rule, which says that $\limsup f/g \ge \liminf f'/g'$ when $\lim f = \lim g = 0$, to the formula (26) for $\pi(x)$ to find

$$\limsup_{x \downarrow 0} \pi(x) \ge \liminf_{x \downarrow 0} \frac{\frac{\mathrm{d}}{\mathrm{d}x} \int_{q_{1-x}}^{\infty} \varphi(G(y) - G(-y)) \,\mathrm{d}G(y)}{\frac{\mathrm{d}}{\mathrm{d}x} \int_{1-x}^{1} \varphi(y) \,\mathrm{d}y}$$
(27)

$$= \liminf_{x \downarrow 0} \frac{\varphi(1-x)g(q_{1-x})}{\varphi(1-x)} \cdot \frac{1}{g(q_{1-x}) + g(-q_{1-x})},$$
(28)

where the equality uses the fundamental theorem of calculus and (25). Condition (P4) on φ implies $\varphi(q)$ must be positive on a neighborhood $q \in (1 - \epsilon, 1)$ for some $\epsilon > 0$, which ensures the limit is well-defined. Reparametrizing in terms of $q = q_{1-x}$, and noting that $q_{1-x} \uparrow \infty$ as $x \downarrow 0$ since g is positive throughout \mathbb{R} , we have

$$\limsup_{x \downarrow 0} \pi(x) \ge \liminf_{q \uparrow \infty} \frac{1}{1 + \frac{g(-q)}{g(q)}}.$$
(29)

Hence $\limsup_{x\downarrow 0} \frac{\pi(x)}{1-\pi(x)} \ge \liminf_{q\uparrow\infty} \frac{g(q)}{g(-q)}$. The conclusion follows from Theorem 2.

Plugging the normal density into Theorem 3 for g(x) confirms the fact suggested by Figure 3:

Corollary 1. If $G = \mathcal{N}(\tau, \sigma^2)$, then $\widetilde{\Gamma}_{\varphi, \text{unif}} = \infty$. That is, no matter what value of Γ is used in a sensitivity analysis with a uniform general signed rank test, the power under $H_1(G)$ tends to one as $n \to \infty$.

5 Simulations

Figures 4 and 5 illustrate Theorem 2 with simulations under standard normal, Laplace and Cauchy alternatives; in each case $\tau = 1/2$, except for the "rare effects" panels in Figure 4 which use the rare effects model described in Section 4. We simulate both standard, fixed-sample tests and uniform tests based on Theorem 1, with the four score functions introduced in Section 2.2. All tests are run with level $\alpha = 0.05$ and plots are based on 10,000 replications.

The results are consistent with our findings above. Figure 4 compares power for each uniform test to the the corresponding fixed-sample test based on the same score function. In the normal case, the uniform test does not indicate finite design sensitivity, as we expect from Corollary 1, and all uniform tests show substantial gains over their fixed-sample counterparts for $n \ge 1,000$. In the Laplace and Cauchy cases, the uniform sign test still shows gains, but uniform tests based on other score functions often fail to outperform their fixed-sample counterparts, as we expect from Figure 3. With large sample sizes, however, the uniform tests at least remain competitive in nearly all cases. Finally, the "rare effects" case again confirms our expectations from Figure 3, showing that each uniform test improves substantially on its fixed-sample counterpart, even with Cauchy noise. Though not shown, the gains for normal and Laplace noise under the rare effects model are even more dramatic, as one would expect by Figure 3.

Figure 5 compares power between uniform tests with different score functions. Tests tend to perform similarly with small sample sizes, but clear distinctions emerge with large sample sizes. In the normal case, the normal scores test dominates while the redescending score function substantially underperforms. As we have seen, under normal noise the outliers contain the most information, and a score function which places more weight upon pairs with large absolute differences will attain higher power as a result. Conversely, in the Cauchy case, the normal scores tests performs the worst, while the sign test performs the best. Here the extreme tails yield less information, as indicated by Figure 3. The Laplace case is a middle ground in which the tails yield no more or less information than most of the rest of the distribution, as we have seen in Figure 3. Here the choice of score function makes little difference.

We close by noting that the uniform sign test shows considerable promise for use in practice. It is competitive in all cases and is the strongest performer of the four tests considered here in a number of cases. This is particularly interesting since the fixed-sample sign test is arguably the least attractive among the fixedsample tests we have considered. It seems the landscape of uniform general signed rank tests is qualitatively different from that of their fixed-sample counterparts.

6 Handling ties

Under the assumption that outcomes are drawn from a continuous distribution, ties among outcome observations occur with probability zero. In practice however, tied outcome data may arise in a variety of settings. In this section we discuss how to adapt the results of the paper to the setting of ties.





Figure 4: Comparison of simulated power for fixed-sample tests (dashed lines) vs. uniform tests (solid lines), based on 10,000 replications. "Cauchy rare effects" panels show "rare effects" alternative model based on Cauchy distribution, as described in Section 4. Other panels show alternative model $H_1(G)$ with distribution G as indicated, having center 1/2 and unit scale. All tests use $\alpha = 0.05$.



—— Sign – – – WSRT …… Normal scores … – – Redescending

Figure 5: Comparison of simulated power for uniform tests using different score functions, based on 10,000 replications under alternative model $H_1(G)$ with G as indicated, having center 1/2 and unit scale. All tests use $\alpha = 0.05$.

Let $Y_{(1)}, \ldots, Y_{(n)}$ be the outcome data ordered in any way so that $|Y_{(1)}| \leq |Y_{(2)}| \leq \ldots |Y_{(n)}|$. Note that this ordering is not unique when ties are present; in such cases, choose one such ordering arbitrarily. We may still apply the methods described in the paper directly to conduct a test. The test statistic and the uniform bound are clearly defined given our chosen ordering of outcomes, and Theorem 1 holds since no aspect of its proof depends on the absence of ties. We remark that it is reasonable to expect $\mathbb{P}(Y_i = 0) > 0$ in the presence of ties; however, this only reduces $\mathbb{P}(Y_i > 0 | \mathcal{F})$, so Proposition 1 and Theorem 1 continue to hold.

However, the version of our uniform test in Theorem 1 depends on the ordering we choose, perhaps arbitrarily, for $(Y_{(i)})$. To remove this undesirable feature of the procedure, we may instead use a generalization of $T_n(x)$ which is invariant to the specific choice of ordering in the tied setting. We write $T_n^*(x)$ for this new test statistic. The intuition for T_n^* comes from recognizing that when ties are present, one or more test statistics in the family $(T_n(x))_{x \in (0,1)}$ are partial sums that include some terms with a particular absolute value but exclude others with the same absolute value, and that the scores associated with these terms may be different from one another. We obtain the family $(T_n^*(x))_{x \in (0,1)}$ by replacing the score for each tied value by the average of scores for all indices involved in the tie, and by adding all these terms together to the partial sum rather than allowing partial sums that contain some terms but not all.

Formally, define $\mathcal{J}_i = \{j \in [n] : |Y_{(j)}| = |Y_{(i)}|\}$, the set of ranks with equal absolute pair differences to the i^{th} ranked pair. Let $m(i) = \min \mathcal{J}_i$, the lowest rank within the tied group containing the i^{th} ranked pair. Now define the test statistic

$$T_{\varphi}^{\star}(x) \coloneqq \sum_{\{i:m(i) \ge (1-x)(n+1)\}} c_i^{\star} \mathbf{1}_{Y_{(i)} > 0}, \quad \text{where} \quad c_i^{\star} = \left|\mathcal{J}_i\right|^{-1} \sum_{j \in \mathcal{J}_i} \varphi\left(\frac{j}{n+1}\right). \tag{30}$$

When a group of pairs share the same absolute outcome value, this test statistic treats all these pairs as a single unit, including either all or none of them in the partial sum, and assigning each a score equal to the average score across all members in the tied set. Note that if there are only k < n distinct absolute outcome

values, there are only k distinct nontrivial values for $T_n^*(x)$; however, if no ties are present, T_n^* is identical to T_n .

We obtain a uniform boundary for $T_n^{\star}(x)$ by substituting c_i^{\star} for c_i in (9) and (10), yielding new quantities $\sigma_n^{\star 2}(x)$ and $f_{\alpha,n}^{\star}(x)$. In the absence of ties, the quantities $\sigma_n^{\star 2}$, and $f_{\alpha,n}^{\star}$ coincide with the original quantities σ_n^{\star} , and $f_{\alpha,n}$. However, the quantities $(c_i^{\star})_{i=1}^n$ are random, unlike (c_i) , hence σ_n^{\star} , and $f_{\alpha,n}^{\star}$ are random as well. This requires no real change to the analysis, since these quantities are \mathcal{F} -measurable and we condition on \mathcal{F} throughout the proof of Theorem 1. As the reader may expect, the new boundary $f_{\alpha,n}^{\star}$ yields a valid uniform test of the sensitivity null $H_0(\Gamma)$ using the order-invariant test statistic T_n^{\star} .

Theorem 4. Under $H_0(\Gamma)$, for any \mathcal{F} -measurable $x_0 > 0$ such that $\sigma_n^{\star 2}(x_0) > 0$ a.s., and any $\alpha \in (0,1)$, we have $\mathbb{P}\left(\exists x \in (0,1) : T_n^{\star}(x) \ge f_{\alpha,n}^{\star}(x) \mid \mathcal{F}\right) \le \alpha$.

Proof. Write $\widetilde{T}_n(x) \coloneqq \sum_{i=\lceil (1-x)(n+1)\rceil}^n c_i^* \mathbf{1}_{Y_{(i)}>0}$; this is the same as $T_n(x)$ with c_i^* substituted for c_i . Repeating the proof of Theorem 1 with σ_n^* and $f_{\alpha,n}^*$ in place of their unstarred counterparts, we obtain

$$\mathbb{P}\left(\left.\exists x \in (0,1) : \widetilde{T}_n(x) \ge f^{\star}_{\alpha,n}(x) \right| \mathcal{F}\right) \le \alpha.$$
(31)

Since $m(i) \leq i$ and $c_i \geq 0$ for all i, we have $T_n^{\star}(x) \leq \tilde{T}_n(x)$ for all x, which implies the result together with (31).

7 Application: impact of fish consumption on mercury concentration

Mercury can be harmful to human health when concentrated too heavily in the bloodstream. There is a substantial body of evidence that consuming large amounts of fish can lead to elevated levels of mercury in the blood (Mahaffey et al., 2004). To study the impact of a high-fish diet on mercury concentration in the blood, we use data from the National Health and Nutrition Examination Survey or NHANES (Centers for Disease Control and Prevention (CDC) National Center for Health Statistics (NCHS), 2017), which records detailed information about respondents' diets and also contains analysis of blood samples, including a measure of total mercury concentration. We identified all 1,672 NHANES respondents from 2007 to 2016 who consumed an average of 15 or more servings of fish monthly, and matched each one to a similar respondent who consumed two or fewer servings of fish per month. Respondents were matched only to respondents from the same two-year period (2007-2008, 2009-2010, etc.). Within these groups, pairs were chosen by optimal matching with respect to a robust Mahalanobis distance (Rosenbaum, 2010a, sec. 8.3) computed from respondent age, household income, gender, ethnicity, cigarettes smoked per day, and indicators for high school graduation, missing high school graduation status, and smoking more than 7 cigarettes per day. Matches were also required to obey a propensity score caliper of 0.2 standard deviations based on a propensity score fitted to these same variables (Rosenbaum and Rubin, 1985). The final matched sample of 1,672 pairs achieved a high degree of balance on covariates, as shown in Table 1. Matching was conducted using R packages rcbalance and optmatch with package cobalt used for balance checking (Pimentel, 2016; Hansen and Klopfer, 2006; Greifer, 2018). For more discussion on the optimal construction of matched samples see Rosenbaum (1989), Hansen (2004), Zubizarreta et al. (2014), and Pimentel et al. (2015).

Note that although the balance on observed variables in Table 1 is very close, individuals with high-fish diets may differ from individuals with low-fish diets on many unobserved attributes correlated with mercury levels. Accordingly, we are interested not only in whether a test assuming an absence of unobserved confounders rejects the null hypothesis, but in how sensitive such a result is to potential bias from unobserved confounders.

In each of the 1,672 pairs formed, we computed the difference in total mercury concentration (in micrograms per mole) between the respondent with the high-fish diet and the respondent with the low-fish diet. The average concentration for matched individuals with high-fish diets and low-fish diets were 3.76 and 1.02 respectively, yielding an average pair difference of 2.73 micrograms per mole. We next tested the sharp null of no effect of treatment in any pair. Mercury measurements were rounded to two decimal places which led to some ties, so for each test we used the test statistic $T_n^*(x)$ of Section 6 and $x_0 = 1/3$ in Theorem 4.

The first three columns of Table 2 show the results of sensitivity analysis in the matched data for the four general signed rank tests considered in this paper. For each of these test statistics, the naïve test with

	Average attr	ibute values	Standardized
Variable	15+ fish servings / mo	0-2 fish servings / mo	difference
Age	43.73	43.63	0.005
Household Income/ $(2x \text{ poverty line})$	2.99	2.96	0.017
Female	0.46	0.46	0.004
Hispanic	0.19	0.18	0.002
Black	0.22	0.22	0.001
Smoker	0.44	0.42	0.011
Cigarettes/Day	4.09	4.04	0.011
High School Graduate	0.80	0.80	0.000
Missing HS Graduation Status	0.03	0.03	0.000

Table 1: Balance table for 1,672 matched pairs formed from NHANES data. Each pair contains one individual who consumed ≥ 15 servings of fish in the previous month, and one who consumed no more than two. The first two columns give the sample means in the matched samples for various attributes of interest, and the third gives the standardized difference, which is computed by dividing the difference in group sample means by the pooled standard deviation estimate from the full dataset before matching.

 $\Gamma = 1$ produces results highly significant at the 0.05 level, and the numbers in the table describe the smallest amount of unmeasured bias necessary to explain the observed effects assuming there is no true effect of treatment—that is, the minimum value of Γ at which we fail to reject the sensitivity analysis null. For example, the fixed-sample sign test ceases to reject the null when we allow for an unobserved confounder which increases the odds of a high-fish diet by a factor of $\Gamma = 4.82$; in contrast, the uniform sign test requires an unobserved confounder which increases the odds of a high fish diet by $\Gamma = 10.51$ before it ceases to reject.

	1,672 pa	airs	190 pairs	
Score function	Fixed-sample	Uniform	Fixed-sample	Uniform
Sign	4.82	10.51	3.72	8.29
Wilcoxon Signed Rank	8.06	10.47	6.04	8.09
Normal Scores	8.55	10.36	6.52	7.95
Redescending	9.68	9.97	7.26	7.58

Table 2: Sensitivity analysis for matched data. Each cell of the table represents a different test statistic for testing the null of no effect of a fish diet on mercury concentration; the first two columns give results for the full matched sample of 1,672 pairs, while the third and fourth columns give results for the smaller sample from 2015-2016 alone. The number in each cell is the smallest degree of unmeasured confounding Γ necessary in the sensitivity analysis model before the test no longer rejects at the $\alpha = 0.05$ level.

Note that repeating the same test many times with different test statistics, as in Table 2, is not recommended in practice. To avoid issues with multiple testing and Type I error control, one should select a single test statistic in advance, possibly based on a pilot sample as described in Heller et al. (2009). We show the results of several tests here to illustrate the impact of the choice of test statistic and complement the discussion in Section 5.

Several interesting patterns are clear in the full-sample results of Table 2. First, regardless of the score function used, the uniform version of the test is less sensitive to unmeasured bias than the fixed-sample version. This pattern is consistent with Theorem 2, which tells us that in large samples the uniform test should perform at least as well as any fixed-sample test it incorporates. Second, the performance of the uniform test across score functions varies much less than the performance of the fixed-sample version across score functions. In particular, the sign test performs substantially worse than any other test examined in the fixed-sample case, but it is comparable to (and even slightly better than) the other score functions in the uniform setting, corroborating the evidence from simulations in Section 5. In this dataset, as in the simulations, adapting over many different truncated statistics appears to compensate for the deficiencies of the fixed-sample sign test.

Finally, we briefly consider the importance of sample size by analyzing the subset of the matched dataset consisting only of those respondents from the final two-year period (2015-2016), a total of 190 pairs. The final two columns of Table 2 repeat the analysis for this smaller dataset. The same pattern of results is observed, with the uniform test outperforming the fixed-sample test for each score function, and the sign

test performing best among uniform tests. Although the benefits of uniform testing articulated in Theorem 2 relate to asymptotic performance in large samples, uniform tests may also offer substantial improvement in datasets of only moderate size.

8 Conclusion and future work

We have described a new test for causal effects in a paired observational study, the uniform general signed rank test. This test provides non-asymptotically valid inference under Rosenbaum's sensitivity analysis model and yields qualitative improvements in design sensitivity relative to existing methods. Our simulation results indicate that the advantages of this test extend from the asymptotic regime down to moderate sample sizes under a variety of alternative hypotheses, as well as to real-world studies.

Though we have described a sensible method for handling ties, we have focused our study on continuous outcomes. When ties are present but rare, as in the data example of Section 7, our findings should continue to hold. However, the study of outcomes with relatively few unique values may require alternative methodology. In such cases, the random walk $(T_n^*(x))_{x \in (0,1)}$ will have relatively few steps, at most the number of unique values of the outcome, with each step comprised of many individual observations, namely all those pairs with absolute outcome equal to a given value. In the sequential analysis literature, such random walks are handled well by group sequential designs (Pocock, 1977; O'Brien and Fleming, 1979; Lan and DeMets, 1983; Jennison and Turnbull, 2000). An application to uniform general signed rank tests may yield promising future results.

Another interesting avenue is the evaluation other theoretical properties, beyond design sensitivity, of uniform general signed rank tests. For example, Lehmann and Romano (2005, Chapter 6) discuss the locally most powerful property of general signed rank tests against particular families of alternatives determined by the function φ . The uniform test is adaptively choosing from among a family of related φ functions, and it would be interesting to understand what the implications are for local optimality in the sense discussed by Lehmann and Romano.

9 Acknowledgments

We thank Eli Ben-Michael for the conversation which sparked this project. Howard thanks Office Of Naval Research (ONR) Grant N00014-15-1-2367.

References

- Boucheron, S., Lugosi, G. and Massart, P. (2013), Concentration inequalities: a nonasymptotic theory of independence, 1st edn, Oxford University Press, Oxford.
- Centers for Disease Control and Prevention (CDC) National Center for Health Statistics (NCHS) (2017), National health and nutrition examination survey data 2009-2016, Technical report, U.S. Department of Health and Human Services, Centers for Disease Control and Prevention, Hyattsville, MD. https: //wwwn.cdc.gov/nchs/nhanes/.
- Chernoff, H. (1952), 'A Measure of Asymptotic Efficiency for Tests of a Hypothesis Based on the sum of Observations', *The Annals of Mathematical Statistics* **23**(4), 493–507.
- Cornfield, J., Haenszel, W., Hammond, E. C., Lilienfeld, A. M., Shimkin, M. B. and Wynder, E. L. (1923/2009), 'Smoking and lung cancer: recent evidence and a discussion of some questions.', *International Journal of Epidemiology* 38(5), 1175–1191.
- Cramér, H. (1938), 'Sur un nouveau théorème-limite de la théorie des probabilités', Actualités Scientifiques 736.

Durrett, R. (2013), Probability: Theory and Examples, 4.1 edn.

- Fogarty, C. B. and Small, D. S. (2016), 'Sensitivity Analysis for Multiple Comparisons in Matched Observational Studies Through Quadratically Constrained Linear Programming', *Journal of the American Statistical Association* 111(516), 1820–1830.
- Gilbert, P. B., Bosch, R. J. and Hudgens, M. G. (2003), 'Sensitivity Analysis for the Assessment of Causal Vaccine Effects on Viral Load in HIV Vaccine Trials', *Biometrics* 59(3), 531–541.

- Greifer, N. (2018), cobalt: Covariate Balance Tables and Plots. R package version 3.5.0, https://CRAN. R-project.org/package=cobalt.
- Hansen, B. B. (2004), 'Full matching in an observational study of coaching for the SAT', Journal of the American Statistical Association 99(467), 609–618.
- Hansen, B. B. and Klopfer, S. O. (2006), 'Optimal full matching and related designs via network flows', Journal of Computational and Graphical Statistics 15(3), 609–627.
- Heller, R., Rosenbaum, P. R. and Small, D. S. (2009), 'Split samples and design sensitivity in observational studies', 104(487), 1090–1101.
- Hewitt, E. and Stromberg, K. R. (1965), Real and Abstract Analysis, Springer-Verlag.
- Howard, S. R., Ramdas, A., McAuliffe, J. and Sekhon, J. (2018a), 'Exponential line-crossing inequalities', arXiv:1808.03204 [math].
- Howard, S. R., Ramdas, A., McAuliffe, J. and Sekhon, J. (2018b), 'Uniform, nonparametric, non-asymptotic confidence sequences', arXiv:1810.08240 [math, stat].
- Jennison, C. and Turnbull, B. W. (2000), Group sequential methods with applications to clinical trials, Chapman & Hall/CRC, Boca Raton.
- Knapp, A. W. (2007), Basic Real Analysis, Springer Science & Business Media.
- Lan, K. K. G. and DeMets, D. L. (1983), 'Discrete Sequential Boundaries for Clinical Trials', *Biometrika* **70**(3), 659–663.
- Lehmann, E. L. and Romano, J. P. (2005), Testing statistical hypotheses, 3rd ed edn, Springer, New York.
- Mahaffey, K. R., Clickner, R. P. and Bodurow, C. C. (2004), 'Blood organic mercury and dietary mercury intake: National health and nutrition examination survey, 1999 and 2000.', *Environmental health* perspectives 112(5), 562.
- Noether, G. E. (1973), 'Some Simple Distribution-Free Confidence Intervals for the Center of a Symmetric Distribution', *Journal of the American Statistical Association* **68**(343), 716–719.
- O'Brien, P. C. and Fleming, T. R. (1979), 'A Multiple Testing Procedure for Clinical Trials', *Biometrics* **35**(3), 549–556.
- Pimentel, S. D. (2016), 'Large, sparse optimal matching with r package rcbalance', *Observational Studies* 2, 4–23.
- Pimentel, S. D., Kelz, R. R., Silber, J. H. and Rosenbaum, P. R. (2015), 'Large, sparse optimal matching with refined covariate balance in an observational study of the health outcomes produced by new surgeons', *Journal of the American Statistical Association* 110(510), 515–527.
- Pocock, S. J. (1977), 'Group Sequential Methods in the Design and Analysis of Clinical Trials', *Biometrika* **64**(2), 191–199.
- Robins, J. M., Rotnitzky, A. and Scharfstein, D. O. (2000), Sensitivity Analysis for Selection bias and unmeasured Confounding in missing Data and Causal inference models, *in* M. E. Halloran and D. Berry, eds, 'Statistical Models in Epidemiology, the Environment, and Clinical Trials', The IMA Volumes in Mathematics and its Applications, Springer New York, pp. 1–94.
- Rosenbaum, P. R. (1989), 'Optimal matching for observational studies', Journal of the American Statistical Association 84(408), 1024–1032.
- Rosenbaum, P. R. (2002), *Observational Studies*, Springer Series in Statistics, 2nd edn, Springer, New York, NY.
- Rosenbaum, P. R. (2004), 'Design Sensitivity in Observational Studies', Biometrika 91(1), 153–164.
- Rosenbaum, P. R. (2010a), *Design of Observational Studies*, Springer Series in Statistics, Springer, New York, NY.
- Rosenbaum, P. R. (2010b), 'Design Sensitivity and Efficiency in Observational Studies', Journal of the American Statistical Association 105(490), 692–702.
- Rosenbaum, P. R. (2011), 'A New u-Statistic with Superior Design Sensitivity in Matched Observational Studies', *Biometrics* 67(3), 1017–1027.
- Rosenbaum, P. R. (2012), 'An exact adaptive test with superior design sensitivity in an observational study of treatments for ovarian cancer', *The Annals of Applied Statistics* 6(1), 83–105.
- Rosenbaum, P. R. and Rubin, D. B. (1985), 'Constructing a control group using multivariate matched sampling methods that incorporate the propensity score', *The American Statistician* **39**(1), 33–38.
- Rosenbaum, P. R. and Small, D. S. (2017), 'An adaptive Mantel-Haenszel test for sensitivity analysis in observational studies', *Biometrics* 73(2), 422–430.
- Sen, P. K. (1970), 'On Some Convergence Properties of One-Sample Rank Order Statistics', The Annals of Mathematical Statistics 41(6), 2140–2143.
- Ville, J. (1939), Étude Critique de la Notion de Collectif., Gauthier-Villars, Paris.

- Yu, B. and Gastwirth, J. L. (2005), 'Sensitivity analysis for trend tests: application to the risk of radiation exposure', *Biostatistics* 6(2), 201–209.
- Zubizarreta, J. R., Paredes, R. D., Rosenbaum, P. R. et al. (2014), 'Matching for balance, pairing for heterogeneity in an observational study of the effectiveness of for-profit and not-for-profit high schools in chile', *The Annals of Applied Statistics* 8(1), 204–231.

A Additional proofs

A.1 Proof of Proposition 1

Throughout the proof, we condition on \mathcal{F} , dropping it from the notation for simplicity. For each $i \in [n]$, write $S_i \coloneqq 1_{Y_{(i)}>0}, X_i \coloneqq T_n(i/(n+1)) = \sum_{j=n-i+1}^n c_j S_j$, and $a_i \coloneqq f_{\alpha,n}(i/(n+1))$. Under $H_0(\Gamma)$, the (S_i) are independent with $1/(1+\Gamma) \leq \mathbb{P}(S_i = 1) \leq \Gamma/(1+\Gamma)$. Let $p_i \coloneqq \mathbb{P}(S_i = 1)$. We wish to show that the rejection probability $\mathbb{P}(\exists i \in [n] : X_i \geq a_i)$ is maximized when $p_i = \Gamma/(1+\Gamma)$ for all $i \in [n]$.

Write $S := (S_1, \ldots, S_n)^T$, a random vector in $\{0, 1\}^n$. Note that, for $s \in \{0, 1\}^n$, $\mathbb{P}(S = s) = \prod_{i=1}^n p_i^{s_i} (1 - p_i)^{1-s_i}$. Let $\mathcal{R} := \left\{s \in \{0, 1\}^n : \sum_{j=n-i+1}^n c_j s_j \ge a_i \text{ for some } i \in [n]\right\}$. This set represents the rejection event, in the sense that the test rejects if and only if $S \in \mathcal{R}$. We will show that $\mathbb{P}(S \in \mathcal{R})$ is increasing in p_i for each $i \in [n]$, from which it follows that the rejection probability is maximized when p_i is maximized for each i.

We claim that if $s \in \mathcal{R}$ and $s' \geq s$ elementwise, then $s' \in \mathcal{R}$. To see this, observe that $s \in \mathcal{R}$ implies that we can choose $i \in [n]$ such that $\sum_{j=n-i+1}^{n} c_j s_j \geq a_i$. Then $\sum_{j=n-i+1}^{n} c_j s_j' \geq \sum_{j=n-i+1}^{n} c_j s_j \geq a_i$, so $s' \in \mathcal{R}$.

Now write $\mathbb{P}(S \in \mathcal{R}) = \sum_{s \in \mathcal{R}} \prod_{i=1}^{n} p_i^{s_i} (1 - p_i)^{1 - s_i}$, and differentiate with respect to p_k for any $k \in [n]$:

$$\frac{\mathrm{d}}{\mathrm{d}p_k} \mathbb{P}(S \in \mathcal{R}) = \sum_{s \in \mathcal{R}} \left[(2s_k - 1) \prod_{i \neq k} p_i^{s_i} (1 - p_i)^{1 - s_i} \right]$$
(32)

$$= \sum_{\substack{s \in \mathcal{R} \\ s_k = 1}} \pi^{(k)}(s) - \sum_{\substack{s \in \mathcal{R} \\ s_k = 0}} \pi^{(k)}(s),$$
(33)

where $\pi^{(k)}(s) = \prod_{i \neq k} p_i^{s_i} (1-p_i)^{1-s_i}$. For each $s \in \mathcal{R}$ with $s_k = 0$, there corresponds an s' which is identical except for $s'_k = 1$, i.e., $s'_i = s_i \mathbf{1}_{i \neq k} + \mathbf{1}_{i=k}$, and this $s' \in \mathcal{R}$ by the claim above. Also, $\pi^{(k)}(s) = \pi^{(k)}(s')$. Hence each term in the second sum of (33) is canceled by a term in the first sum. We conclude $\frac{d}{dp_k} \mathbb{P}(S \in \mathcal{R}) \ge 0$, as desired.

We remark that an alternative proof could use Holley's inequality for distributions over finite distributive lattices (Rosenbaum, 2002, Sections 2.10, 4.7.2). We have opted for the direct proof above to keep the paper more self-contained.

A.2 A technical result on Riemann sums

The following result ensures convergence of certain Riemann sums for some unbounded functions, and is necessary to analyze the asymptotic behavior of $f_{\alpha,n}$.

Lemma 3. Suppose $\varphi : (0,1) \to [0,\infty)$ is discontinuous on a set of measure zero, $\int_0^1 \varphi(x) dx < \infty$, and there exists a constant $a \in [0,1/2)$ such that φ is nonincreasing on (0,a), nondecreasing on (1-a,1), and bounded on (a, 1-a). Then

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \varphi\left(\frac{i}{n+1}\right) = \int_{0}^{1} \varphi(x) \,\mathrm{d}x.$$
(34)

Proof. Write $\varphi = \varphi_1 + \varphi_2 + \varphi_3$ where $\varphi_1(x) \coloneqq \varphi(x) \mathbf{1}_{x < a}$, $\varphi_2(x) \coloneqq \varphi(x) \mathbf{1}_{a \le x \le 1-a}$, and $\varphi_3(x) \coloneqq \varphi(x) \mathbf{1}_{x > a}$. Since φ_2 is bounded, it is Riemann integrable, so $n^{-1} \sum_{i=1}^n \varphi_2(i/(n+1)) \to \int_0^1 \varphi_2(x) \, \mathrm{d}x$ by standard Riemann integration theory, noting that $i/(n+1) \in ((i-1)/n, i/n)$ for each $i \in [n]$. For φ_1 and φ_3 , we appeal to Lemma 4 below to conclude that $n^{-1} \sum_{i=1}^{n} \varphi_k(i/(n+1)) \to \int_0^1 \varphi_k(x) \, dx$ for k = 1, 3. The result follows by linearity.

Lemma 4. Suppose $\varphi: (0,1) \to [0,\infty)$ is monotone and $\int_0^1 \varphi(x) dx < \infty$. Then

$$\lim_{n \to \infty} \frac{1}{n} \sum_{i=1}^{n} \varphi\left(\frac{i}{n+1}\right) = \int_{0}^{1} \varphi(x) \,\mathrm{d}x.$$
(35)

Proof. Suppose first that φ is nondecreasing, and for each $n \in \mathbb{N}$ define $\varphi_n(x) \coloneqq \varphi(i/(n+1))$ for $i/(n+1) \leq x < (i+1)/(n+1)$, i = 1, ..., n, and $\varphi_n(x) = 0$ for x < 1/(n+1). Then $|\varphi_n| \leq |\varphi|$ for all n by construction, since φ is nonnegative and nondecreasing. Furthermore, since φ is monotone, it is discontinuous at a countable number of points (Knapp, 2007, p. 344), so $\varphi_n(x) \to \varphi(x)$ pointwise almost everywhere. So the dominated convergence theorem implies

$$\lim_{n \to \infty} \frac{1}{n+1} \sum_{i=1}^{n} \varphi\left(\frac{i}{n+1}\right) = \lim_{n \to \infty} \int_{0}^{1} \varphi_n(x) \,\mathrm{d}x = \int_{0}^{1} \varphi(x) \,\mathrm{d}x.$$
(36)

The conclusion follows since $(n+1)/n \to 1$ as $n \to \infty$. If φ is instead nonincreasing, apply the above argument to $x \mapsto \varphi(1-x)$.

A.3 Proof of Lemma 1

The limit $n^{-1}\mu_n(x) \to \rho_{\Gamma} \int_{1-x}^1 \varphi(y) \, dy$ follows directly from Lemma 3 applied to the function $q \mapsto \varphi(1-q) \mathbf{1}_{q \leq x}$. The bulk of the work is in proving that $f_{\alpha,n}(x) = \mu_n(x) + \mathcal{O}(\sqrt{n})$. For this, fix $\rho \in [1/2, 1)$ and let $h(x) := e^x / [1 + \rho(e^x - 1)]^2$. We require the following technical lemma, proved below.

Lemma 5. For any $\rho \in [1/2, 1)$, $0 \le h(x) \le 1$ for all $x \ge 0$.

To prove Lemma 1, we use a first-order application of Taylor's theorem about $\lambda = 0$, which yields, for any $c \ge 0, \lambda \ge 0$,

$$\log\left(1+\rho(e^{c\lambda}-1)\right) = \rho c\lambda + \frac{\rho(1-\rho)h(\xi)c^2\lambda^2}{2},\tag{37}$$

for some $\xi \in [0, c\lambda]$. Since $\Gamma \ge 1$, we are assured $\rho_{\Gamma} \ge 1/2$, as assumed above. So combining the definition (10) of $f_{\alpha,n}$ with the expansion (37), we have

$$f_{\alpha,n}(x) = \frac{\log \alpha^{-1}}{\lambda_n} + \mu_n(x) + \frac{\rho_{\Gamma}(1-\rho_{\Gamma})\lambda_n}{2} \sum_{i=\lceil (1-x)(n+1)\rceil}^n h(\xi_i)c_i^2,$$
(38)

where $\xi_i \in [0, c_i \lambda_n]$ for each i = 1, ..., n. Now Lemma 5 implies

$$0 \le \frac{\rho_{\Gamma}(1-\rho_{\Gamma})\lambda_n}{2} \sum_{i=\lceil (1-x)(n+1)\rceil}^n h(\xi_i)c_i^2 \le \frac{\lambda_n \sigma_n^2(x)}{2},\tag{39}$$

so that

$$0 \le f_{\alpha,n}(x) - \mu_n(x) \le \frac{\log \alpha^{-1}}{\lambda_n} + \frac{\lambda_n \sigma_n^2(x)}{2}.$$
(40)

Applying Lemma 3 to the function $q \mapsto \varphi^2(1-q)\mathbf{1}_{x \leq x}$, which is integrable by (P1), we see that $n^{-1}\sigma_n^2(x) = \mathcal{O}(1)$ for each $x \in (0, 1)$. Together with the definition (10) of λ_n , we conclude

$$0 \le \frac{f_{\alpha,n}(x) - \mu_n(x)}{\sqrt{n}} = \frac{1}{\sqrt{n}} \left(\frac{\log \alpha^{-1}}{\lambda_n} + \frac{\lambda_n \sigma_n^2(x)}{2} \right) = \sqrt{\frac{\sigma_n^2(x_0)}{2n}} + \sqrt{\frac{2n \log \alpha^{-1}}{\sigma_n^2(x_0)}} \cdot \frac{\sigma_n^2(x)}{n} = \mathcal{O}(1), \quad (41)$$

as desired.

Note that, if we further assume $\int_0^1 \varphi^3(x) dx < \infty$, then we have the second-order expansion mentioned in Section 3,

$$f_{\alpha,n}(x) = \mu_n(x) + \left(1 + \frac{\sigma_n^2(x)}{\sigma_n^2(x_0)}\right) \sqrt{\frac{\sigma_n^2(x_0)\log\alpha^{-1}}{2}} + \mathcal{O}(1).$$
(42)

To prove (42) we follow an analogous argument starting from

$$\log\left(1 + \rho(e^{c\lambda} - 1)\right) = \rho c\lambda + \frac{\rho(1 - \rho)c^2\lambda^2}{2} - \frac{\rho(1 - \rho)h_2(\xi)c^3\lambda^3}{6},\tag{43}$$

for some $\xi \in [0, c\lambda]$, where

$$h_2(x) \coloneqq \frac{e^x [\rho(1+e^x) - 1]}{[1+\rho(e^x - 1)]^3} \tag{44}$$

satisfies $0 \le h_2(x) \le 1$ for all $x \ge 0$. By the same argument which led from (37) to (41), we find

$$0 \le \frac{\log \alpha^{-1}}{\lambda_n} + \mu_n(x) + \frac{\lambda_n \sigma_n^2(x)}{2} - f_{\alpha,n}(x) \le \frac{\rho_{\Gamma}(1-\rho_{\Gamma})\lambda_n^2}{6} \sum_{i=\lceil (1-x)(n+1)\rceil}^n c_i^3 = \mathcal{O}(1).$$
(45)

Substituting the definition of λ_n shows that

$$\frac{\log \alpha^{-1}}{\lambda_n} + \mu_n(x) + \frac{\lambda_n \sigma_n^2(x)}{2} = \mu_n(x) + \left(1 + \frac{\sigma_n^2(x)}{\sigma_n^2(x_0)}\right) \sqrt{\frac{\sigma_n^2(x_0)\log \alpha^{-1}}{2}} \rightleftharpoons g_{\alpha,n}(x).$$
(46)

Note that the chosen value of λ_n is the minimizer of the left-hand side of (46) when $x = x_0$, justifying the claim that λ_n is chosen to optimize the bound $g_{\alpha,n}(x)$ at $x = x_0$.

Proof of Lemma 5. That $h(x) \ge 0$ for all $x \ge 0$ is clear from the definition. To see that $h(x) \le 1$, observe

$$h'(x) = e^x \left(\frac{1 - \rho(1 + e^x)}{[1 + \rho(e^x - 1)]^3} \right).$$
(47)

Now the inequality $e^x \ge 1 + x$ implies $1 - \rho(1 + e^x) \le 1 - 2\rho \le 0$ by our assumption $\rho \ge 1/2$, while $1 + \rho(e^x - 1) \ge 1 > 0$. Hence $h'(x) \le 0$ for all $x \ge 0$. Together with h(0) = 1, the conclusion follows. \Box

A.4 Proof of Lemma 2

Let H(x) := G(x) - G(-x). Fix any $\epsilon > 0$. Because bounded, continuous functions with compact support are dense in L^p (Hewitt and Stromberg, 1965, Theorem 13.21), we can find a continuous function $\varphi_c : [0, 1] \rightarrow [0, \infty)$ such that $\int_0^1 |\varphi(x) - \varphi_c(x)| \, dx < \epsilon$, and $\varphi_c(x) = 0$ for all $x \in [0, b) \cup (1 - b, 1]$ for some 0 < b < a. Now write

$$\tau \coloneqq \int_0^\infty \varphi(H(x)) \,\mathrm{d}G(x) \quad \text{and} \tag{48}$$

$$\tau_c := \int_0^\infty \varphi_c(H(x)) \,\mathrm{d}G(x). \tag{49}$$

We will show

$$\limsup_{n \to \infty} \frac{1}{n} \left| \sum_{i=1}^{n} \varphi\left(\frac{i}{n+1}\right) \mathbf{1}_{Y_{(i)} > 0} - \sum_{i=1}^{n} \varphi_c\left(\frac{i}{n+1}\right) \mathbf{1}_{Y_{(i)} > 0} \right| < \epsilon, \quad \text{a.s.}, \tag{50}$$

$$\frac{1}{n} \sum_{i=1}^{n} \varphi_c \left(\frac{i}{n+1} \right) \mathbf{1}_{Y_{(i)} > 0} \xrightarrow{\text{a.s.}} \tau_c, \quad \text{and}$$
(51)

$$|\tau_c - \tau| < \epsilon, \tag{52}$$

from which we conclude $\limsup_{n\to\infty} \left| n^{-1} \sum_{i=1}^{n} \varphi\left(\frac{i}{n+1}\right) \mathbf{1}_{Y_{(i)}>0} - \tau \right| < 2\epsilon$ a.s. Since ϵ was arbitrary, the conclusion follows.

To obtain (50), use the triangle inequality to write

$$\frac{1}{n} \left| \sum_{i=1}^{n} \left[\varphi\left(\frac{i}{n+1}\right) - \sum_{i=1}^{n} \varphi_c\left(\frac{i}{n+1}\right) \right] \mathbf{1}_{Y_{(i)} > 0} \right| \le \frac{1}{n} \sum_{i=1}^{n} \left| \varphi - \varphi_c \right| \left(\frac{i}{n+1}\right)$$
(53)

$$=\frac{1}{n}\sum_{i=1}^{n+1}|\varphi-\varphi_c|\left(\frac{i}{n+1}\right)-\frac{|\varphi-\varphi_c|(1)}{n}\tag{54}$$

$$\rightarrow \int_{0}^{1} \left| \varphi - \varphi_{c} \right| (x) \, \mathrm{d}x < \epsilon, \tag{55}$$

where the limit uses Lemma 3, noting that $|\varphi - \varphi_c|$ is bounded on [b, 1 - b] and monotone elsewhere, and final inequality follows from our choice of φ_c .

The second step (51) follows from Theorem 1 of Sen (1970) applied to φ_c , which we partially restate. See Appendix A.6 for an explanation of why our statement differs from Sen's.

Lemma 6 (Sen, 1970, Theorem 1). Suppose $\varphi_c \in L^1(0, 1)$ is bounded and continuous, and suppose Y_1, Y_2, \ldots are drawn *i.i.d.* from a continuous distribution G. Then

$$\frac{1}{n}\sum_{i=1}^{n}\varphi_{c}\left(\frac{i}{n+1}\right)\mathbf{1}_{Y_{(i)}\geq0} \stackrel{a.s.}{\to} \int_{0}^{\infty}\varphi_{c}(H(x))\,\mathrm{d}G(x).$$
(56)

Finally, to see (52), use the triangle inequality to write

$$\tau_c - \tau | \le \int_0^\infty |\varphi_c - \varphi| \left(H(y) \right) \mathrm{d}G(y) \tag{57}$$

$$\leq \int_0^\infty |\varphi_c - \varphi| (H(y)) \, \mathrm{d}H(y), \tag{58}$$

since $H'(y) = G'(y) + G'(-y) \ge G'(y)$ and the integrand is nonnegative. From this we conclude

$$|\tau_c - \tau| \le \int_0^1 |\varphi_c - \varphi|(u) \, \mathrm{d}u < \epsilon,$$
(59)

by our choice of φ_c .

A.5 Proof of Proposition 2

Fix any $p \ge 1$. A standard Cramér-Chernoff tail bound for the normal distribution (Boucheron et al., 2013, Section 2.2) gives $1 - \Phi(x) \le e^{-x^2/2}$, which implies $\Phi^{-1}(q) \le \sqrt{2\log(1-q)^{-1}}$. Hence

$$\int_{0}^{1} |\varphi(q)|^{p} \,\mathrm{d}q \le 2^{p/2} \int_{0}^{1} [\log(2/(1-q))]^{p/2} \,\mathrm{d}q \tag{60}$$

$$=2^{1+p/2} \int_{\log 2}^{\infty} y^{p/2} e^{-y} \,\mathrm{d}y \tag{61}$$

using the substitution $y = \log(2/(1-q))$. The final integral is upper bounded by $\Gamma(1+p/2)$, using the definition of the Gamma function and non-negativity of the integrand, which completes the proof.

A.6 Discussion of Theorem 1 from Sen (1970)

Sen (1970) assumes only that $\varphi \in L^1(0,1)$ is continuous. Denoting $\varphi_n(x) \coloneqq \varphi(i/(n+1))$ for $(i-1)/n < x \le i/n, i = 1, \ldots, n$, their proof (p. 2141) claims that

$$\lim_{n \to \infty} \int_0^1 |\varphi_n(x)| \, \mathrm{d}x = \int_0^1 |\varphi(x)| \, \mathrm{d}x.$$
(62)



Figure 6: $\pi(x)$ from Theorem 2 for additional score functions not included in Figure 3, when G is standard normal, Laplace (double exponential) or Cauchy, and $\tau = 1/2$.

The conclusion (62) is not true for all continuous $\varphi \in L^1(0, 1)$, as the counterexample below shows. However, noting that $\int_0^1 \varphi_n(x) dx = n^{-1} \sum_{i=1}^n \varphi(i/(n+1))$, our Lemma 3 shows that (62) is true under stronger conditions, and in particular is true for bounded φ . This is the reason we require boundedness in our restatement of Sen's result, Lemma 6.

Let $\varphi(x) = n$ for $1/(n+1) \le x \le 1/(n+1) + 1/(n2^{n+1})$, $n \in \mathbb{N}$. Then $\int_0^1 \varphi(x) \, dx = \sum_{n=1}^{\infty} 2^{-n-1} = 1/2$, hence $\varphi(x) \in L^1$. But $n^{-1} \sum_{i=1}^n \varphi(i/(n+1)) \ge n^{-1} \varphi(1/(n+1)) = 1$ for all n, so $\liminf_{n \to \infty} n^{-1} \sum_{i=1}^n \varphi(i/(n+1)) \ge 1 > 1/2 = \int_0^1 \varphi(x) \, dx$, showing (62) does not hold. This φ is not continuous, but may be replaced with a continuous approximation by standard arguments.

A.7 Additional plots of $\pi(x)$

Figure 6 plots $\pi(x)$ as defined in Theorem 2 for additional score functions not included in Figure 3.